

## LEC 20 MATH 3A3

we usually don't know when  $b_j$  is a causal cause or not.

However, at the very least, it is correlational.

What is the interpretation of a correlational estimate?

Warning: it's weak! Here it is for OLS,  $x_j = \text{blood sugar (mg/dL)}$ ,  $y = \text{blood triglyceride (mmol/L)}$  for  $b_j = 1.17$ ,  $s_{b_j} = 0.25$

"When comparing two observations (A) and (B) which are sampled in the same fashion as the observations in ① were sampled when (A) has  $\frac{\text{blood sugar}}{\text{value of } x_j}$   $\frac{1 \text{ mg/dL}}{\text{unit}}$

larger than (B) but share the same measurement values

otherwise, then (A) is predicted to have an estimated

$\frac{\text{blood triglyceride level}}{\text{value of response } y}$   $\frac{1.17 \pm 0.25 \text{ mmol/L}}{\text{value of } b_j \pm s_{b_j} \text{ unit of } y}$   $\frac{\text{higher}}{\text{if } b_j > 0 \Rightarrow \text{higher or lower}}$

than (B)'s  $\frac{\text{blood triglyceride level}}{\text{value of response } y}$

library in the p covariance.

assuming the blood triglyceride level is  
the response function against

do logistic regression, neural network, Cox ph.

For prediction only  $\hat{y} = g(\mathbf{x})$ , we don't

Care at all about causation! we just

want the features to have non-spurious

correlations with  $y$ ! This is why this topic was

excluded from 3A2. Because it's irrelevant!!

Why is the causal relationship more powerful than the correlational relationship? Causal relationship means if you do the manipulation, you get the results. You don't get this with a correlational relationship!! Again for prediction, it doesn't matter!

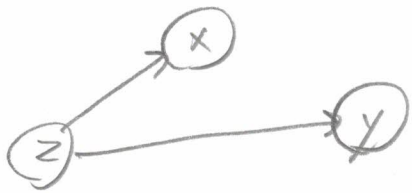
Prediction answers question: what will happen? "You will make more money if  $x$  increases naturally"

Causality answers question: what will make it happen?

"You will make more money if you increase  $x$ ."

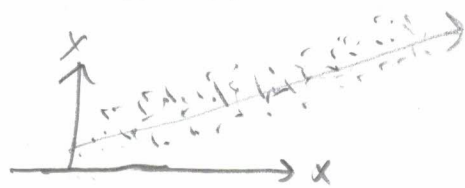
Causality is big business because companies want to know what they can do to make something happen.

When/why could do conditions fail to be causal. Take, [B]



$X = \#$  of umbrellas sold 6-9 AM in NYC  
 $Y = \#$  of car accidents 9 AM-9 PM in NYC  
 $Z = \#$  inches rain forecast in NYC at 6 AM

Obviously  $\text{Corr}[X, Y] \neq 0$  and positive



$b_1$  would be positive and significant

You cannot use this scatterplot as a means to assess causality, you must have other information! You must have the DAG and it must be true!

Let's say the <sup>gods</sup> forecast is and manipulate umbrella sales. They say, you cannot sell to NYC residents, you only sell to me so I can ship them to Paraguay. Then they proceed to give them varying volume orders every day. By manipulating  $X$ , it is now completely disconnected from arrows coming in:



Thus rainfall has no effect on  $X$ , hence there is no information about  $Y$  in  $X$ .



$b_1 \approx 0$  and insignificant

This constitutes proof  $X$  is not causally related to  $Y$ ! We will return to manipulation later!

So far, we've only been estimating single regression.

What happens if we run the regression and fit  $\hat{y} = b_0 + b_1 x + b_2 z$ ?

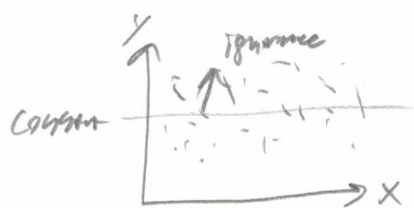
What is interpretation of  $b_1$ ? "but share the same measurement values otherwise ...". This means consider two observations A, B

$x_A = x_0$ ,  $x_B = x_0 + 1$  where  $z_A = z_0$  and  $z_B = z_0$ , what is the difference in the expected response?

Linear regression is a way of looking at a conditional universe!

Consider  $z_A = 17$  and  $z_B = 17$ . The world of all observations

$\langle x, 17, y \rangle$  looks like



$b_1 \approx 0$  and insignificant! Why? We're conditional on rainfall. So  $y = K_{y,z}(17) + \text{ignore}$   
 $= \text{constant} + \text{ignore}$

This is what  $b_1$  in the regression  $\hat{y} = b_0 + b_1 x + b_2 z$  is estimating,  $\beta_1$

It's completely different from  $\beta_1$  in  $\hat{y} = b_0 + b_1 x$ . DEMO

It's a change in perspective, not name itself. Manipulation: change in name itself.

This seems unfair! If we add a feature to a regression can exposure change? Go from positive and significant to zero and insignificant?

Yes! They can do anything. It's all based on the underlying

structure given by the DAG.

# Paradoxes!

You may be thinking, this umbrella-cor-relates with lurking variable rain is obvious! I'll never fall into the trap the

condition  $\Rightarrow$  causation. Not so fast. Consider

$x_1$  = exercise amount and  $y$  = cholesterol level. The data shows:



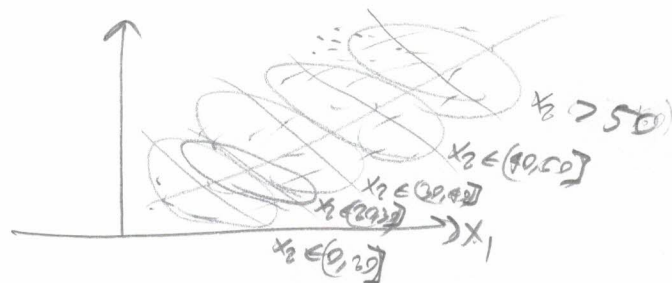
$$\hat{y} = b_0 + b_1 x_1 \Rightarrow$$

$b_1 > 0$  and significant

the data is real!

You're tempted to say more exercise causes higher cholesterol.

How did this happen? There's another variable,  $x_2$  = age.

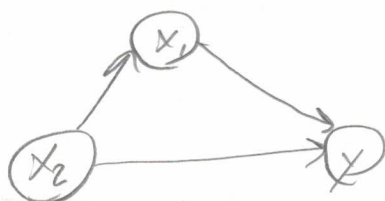


This is called  
"Simpson's  
Paradox"  
(1951)

If you condition on  $x_2$  i.e. run  $\hat{y} = b_0 + b_1 x_1 + b_2 x_2$ ,

then you reveal the real causative effect! PAB? Never...

[D]



$\exists$  a formula

$\Rightarrow$  The effect of  $x_1$  on  $y$  is causal if, direct connection and if there is additionally a common cause, it's controlled for. (You need Multivariate Regression in the latter case.)